

Real-time Visual Tracking Using Sparse Representation

Hanxi Li, Chunhua Shen, and Qinfeng Shi

Abstract—The ℓ_1 tracker obtains robustness by seeking a sparse representation of the tracking object via ℓ_1 norm minimization [1]. However, the high computational complexity involved in the ℓ_1 tracker restricts its further applications in real time processing scenario. Hence we propose a Real Time Compressed Sensing Tracking (RTCST) by exploiting the signal recovery power of Compressed Sensing (CS). Dimensionality reduction and a customized Orthogonal Matching Pursuit (OMP) algorithm are adopted to accelerate the CS tracking. As a result, our algorithm achieves a real-time speed that is up to 6,000 times faster than that of the ℓ_1 tracker. Meanwhile, RTCST still produces competitive (sometimes even superior) tracking accuracy comparing to the existing ℓ_1 tracker. Furthermore, for a stationary camera, a further refined tracker is designed by integrating a CS-based background model (CSBM). This CSBM-equipped tracker coined as RTCST-B, outperforms most state-of-the-arts with respect to both accuracy and robustness. Finally, our experimental results on various video sequences, which are verified by a new metric—Tracking Success Probability (TSP), show the excellence of the proposed algorithms.

Index Terms—Visual tracking, compressed sensing, particle filter, linear programming, hash kernel, orthogonal matching pursuit.

I. INTRODUCTION

Within Bayesian filter framework, the representation of the likelihood model is essential. In a tracking algorithm, the scheme of object representation determines how the concerned target is represented and how the representation is updated. A promising representation scheme should accommodate noises, occlusions and illumination changes in various scenarios. In the literature, a few representation models have been proposed to ease these difficulties [2–7]. Most tracking algorithms represent the target by a single model, typically built on extracted features such as color histogram [8, 9], textures [10] and correspondence points [11]. Nonetheless, these approaches are usually sensitive to variations in target appearance and illumination, and a powerful template update method is usually needed for robustness. Other tracking algorithms train a classifier off-line [5, 12] or on-line [7] based on multiple target samples. These algorithms benefit from the robust object model, which is learned from labeled data by sophisticated learning methods.

H. Li and C. Shen is with NICTA, Canberra Research Laboratory, Canberra, ACT 2601, Australia, and also with the Australian National University, Canberra, ACT 0200, Australia (e-mail: {hanxi.li, chunhua.shen}@nicta.com.au).

Q. Shi is with University of Adelaide, Adelaide, SA 5000, Australia (e-mail: qinfeng.shi@ieee.org).

Correspondence should be addressed to C. Shen.

NICTA is funded by the Australian Government as represented by the Department of Broadband, Communications and the Digital Economy and the Australian Research Council through the ICT Center of Excellence program.

Recently, Mei and Ling proposed a robust tracking algorithm using ℓ_1 minimization [1]. Their algorithm, referred to as the ℓ_1 tracker, is designed within Particle Filter (PF) framework [13]. There a target is expressed as a *sparse* representation of multiple predefined templates. The ℓ_1 tracker demonstrates promising robustness compared with existing trackers [14–16]. However, it has following problems: Firstly, ℓ_1 minimization in their work is slow; Secondly, they use an over-complete dictionary (an identity matrix) to represent the background and noise. This dictionary, in fact, can also represent any objects (including the user interested tracking objects) in video. Hence it may not discriminate the objects against background and noise.

Although the ℓ_1 tracker [1] is inspired by the face recognition work using *sparse representation classification* (SRC)[17], it doesn't make use of the sparse signal recovery power of Compressed Sensing (CS) used in [17]. CS is an emerging topic originally proposed in signal processing community [18, 19]. It states that sparse signals can be exactly recovered with fewer measurements than what the Nyquist-Shannon criterion requires with overwhelming probability. It has been applied to various computer vision tasks [17, 20, 21].

Inspired by the ℓ_1 tracker and motivated by their problems, we propose two CS-based algorithms termed *Real-Time Compressed Sensing Tracking* (RTCST) and *Real-Time Compressed Sensing Tracking with Background Model* (RTCST-B) respectively. The new tracking algorithms are tremendously faster than the standard ℓ_1 tracker and serve as *better* (in terms of both accuracy and robustness) alternatives to existing visual object trackers such as those in [7, 13, 14].

The key contributions of this work can be summarized as follows.

- 1) We make use of the sparse signal recovery power of CS to reduce the computational complexity significantly. That is we hash or random project the original features to a much lower dimensional space to accelerate the CS signal recovery procedure for tracking. Moreover, we propose a customized *Orthogonal Matching Pursuit* (OMP) algorithm for real-time tracking. Our algorithms are up to about 6,000 times faster than the standard ℓ_1 tracker of [1]. In short, *we make the tracker real-time by using CS*.
- 2) We propose background template rather than the over-complete dictionary in [1]. This further improves the robustness of the tracking, because the representation of the objects and background are better separated. This new tracker, which is referred to as RTCST-B in this work, outperforms most state-of-the-art visual trackers

with respect to accuracy while achieves even higher efficiency compared with RTCST.

- 3) Finally, we propose a new metric called *Tracking Success Probability* (TSP) to evaluate trackers' performance. We argue that this new metric is able to measure tracking results quantitatively and demonstrate the robustness of a tracker. Consequently, all the empirical results are assessed by using TSP in this work.

For ease of exposition, symbols and their denotations used in this paper are summarized in Table I.

The rest of the paper is organized as follows. We briefly review the related literature background in the next section. In Section III, the proposed RTCST algorithm is presented. We present the RTCST-B tracker in Section IV. We verify our methods by comparing them against existing visual tracking methods in Section V. Conclusion and discussion can be found in the last section.

II. RELATED WORK

In this section, we briefly review theories and algorithms closest to our work.

A. Bayesian Tracking and Particle Filters

From a Bayesian perspective, the tracking problem is to calculate the posterior probability $p(\mathbf{s}_k|\mathbf{y}_k)$ of state \mathbf{s}_k at time k , where \mathbf{y}_k is the observed measurement at time k [13]. In principle, the posterior PDF is obtained recursively via two stages: prediction and update. The prediction stage involves the calculation of prior PDF:

$$p(\mathbf{s}_k|\mathbf{y}_{k-1}) = \int p(\mathbf{s}_k|\mathbf{s}_{k-1})p(\mathbf{s}_{k-1}|\mathbf{y}_{k-1})d\mathbf{s}_{k-1}. \quad (1)$$

In the update stage, the prior is updated using Bayes' rule

$$p(\mathbf{s}_k|\mathbf{y}_k) = \frac{p(\mathbf{y}_k|\mathbf{s}_k)p(\mathbf{s}_k|\mathbf{y}_{k-1})}{p(\mathbf{y}_k|\mathbf{y}_{k-1})}. \quad (2)$$

The recurrence relations (1) and (2) form the basis for the optimal Bayesian solution. Nonetheless, the solution of above problem can not be analytically solved without further simplification or approximation. Particle Filter (PF) is a Bayesian sequential importance sampling technique for estimating the posterior distribution $p(\mathbf{s}_k|\mathbf{y}_k)$. By introducing the so-called *importance sampling distribution* [8]:

$$\mathbf{s}_i \sim q(\mathbf{s}), \quad i = 1, \dots, N_s, \quad (3)$$

the posterior density is estimated by a weighted approximation,

$$p(\mathbf{s}_k|\mathbf{y}_k) \approx \sum_{i=1}^{N_s} w_k^i \delta(\mathbf{s}_k - \mathbf{s}_k^i). \quad (4)$$

Here

$$w_k^i \propto w_{k-1}^i \frac{p(\mathbf{y}_k|\mathbf{s}_k^i)p(\mathbf{s}_k^i|\mathbf{s}_{k-1}^i)}{q(\mathbf{s}_k^i|\mathbf{s}_{k-1}^i, \mathbf{y}_k)}. \quad (5)$$

For the sake of convenience, $q(\cdot)$ is commonly formed as

$$q(\mathbf{s}_k|\mathbf{s}_{k-1}^i, \mathbf{y}_k) = p(\mathbf{s}_k|\mathbf{s}_{k-1}^i). \quad (6)$$

Therefore, (5) is simplified into

$$w_k^i \propto w_{k-1}^i p(\mathbf{y}_k|\mathbf{s}_k^i) \quad (7)$$

The posterior then could be updated only depending on its previous value and observation likelihood $p(\mathbf{s}_k|\mathbf{s}_{k-1}^i)$. Plus, in order to reducing the effect of *particle degeneracy* [8], a resampling scheme is usually implemented as

$$Pr(\mathbf{s}_k^{i*} = \mathbf{s}_k^j) = w_k^j, \quad j = 1, 2, \dots, N_s \quad (8)$$

where the set $\{\mathbf{s}_k^{i*}\}_{i=1}^{N_s}$ is the particles after re-sampling.

Like the ℓ_1 tracker, both RTCST and RTCST-B trackers use PF framework. However, they differ in how to seek a sparse representation which consequently lead to different observation likelihood $p(\mathbf{s}_k|\mathbf{s}_{k-1}^i)$ estimation.

B. ℓ_1 -norm Minimization-based Tracking

The underlying conception behind SRC is that in many circumstances, an observation belonging to a certain class lies in the subspace that is spanned by the samples belong to this class, and the linear representation is assumed to be sparse. Hence, reconstructing the sparse coefficients associated with the representation is crucial to identify the observation. The coefficients recovery could be accomplished by solving a relaxed version of (13)

$$\min_{\mathbf{x}} \|\mathbf{x}\|_1, \text{ s.t. } \|A\mathbf{x} - \mathbf{y}\|_2 \leq \varepsilon, \quad (9)$$

where $\mathbf{x} \in \mathbb{R}^n$ is the coefficient vector of interest; $A = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n] \in \mathbb{R}^{d \times n}$ is sometimes dubbed as *dictionary* and composed of pre-obtained pattern samples $\mathbf{a}_i \in \mathbb{R}^d \forall i$; and $\mathbf{y} \in \mathbb{R}^d$ is the query/test observation. ε is error tolerance. Then, the class identity $l(\mathbf{y})$ is retrieved as

$$l(\mathbf{y}) = \underset{j \in \{1, \dots, C\}}{\operatorname{argmin}} r_j(\mathbf{y}), \quad (10)$$

where $r_j(\mathbf{y}) \doteq \|\mathbf{y} - A\delta_j(\mathbf{x})\|_2$ is the reconstruction residual associated with class i , C is the number of classes and the function $\delta_j(\mathbf{x})$ sets all the coefficients of \mathbf{x} to 0 except those corresponding to j th class [17].

Given a target template set $T = [\mathbf{t}_1, \dots, \mathbf{t}_{N_t}] \in \mathbb{R}^{d_0 \times N_t}$ and a noise template set $E = [I, -I] \in \mathbb{R}^{d_0 \times 2d_0}$, the ℓ_1 tracker adopts a positive-restricted version of (14) for recovering the sparse coefficients \mathbf{x} , i.e.,

$$\min \|\mathbf{x}\|_1, \text{ s.t. } \|A\mathbf{x} - \mathbf{y}\|_2 \leq \varepsilon, \quad \mathbf{x} \succeq 0. \quad (11)$$

Here $A \doteq [T, E] \in \mathbb{R}^{d_0 \times (N_t + 2d_0)}$ is the combination of target templates and noise templates while $\mathbf{x} \doteq [\mathbf{x}_t^T, \mathbf{x}_e^T]^T \in \mathbb{R}^{N_t + 2d_0}$ denotes the associated target coefficients and noise coefficients. Note that N_t denotes the number of target templates and d_0 is the original dimensionality of feature space which equals to the pixel number of the initial target. The ℓ_1 tracker tracks the target by integrating (11) and a template-update strategy into the PF framework. Algorithm 1 illustrates the tracking procedure. In addition, there is a heuristic approach for updating the target templates and their weights in the ℓ_1 tracker. Refer to [1] for more details.

TABLE I: Notation

Notation	Description
\mathbf{s}_k	A dynamic state vector at time k
\mathbf{s}_k^i	A dynamic state vector at time k corresponding to the i th particle
A	The measurement matrix or the collection of templates
\mathbf{y}	The observed target, a.k.a, observation
\mathbf{x}	The signal to be recovered in compressed sensing. For CS-based pattern recognition or tracking, it is the coefficient vector for the sparse representation
Φ	The projection matrix, could be either a random matrix or a hash matrix in this work
T, E, B	The collection of target, noise and background templates
$\mathbf{x}_t, \mathbf{x}_e, \mathbf{x}_b$	The coefficient vector associated with target, noise and background templates respectively
N_t, N_b	The number of target templates and background templates
d_0, d	The dimensionality of original and reduced feature space

Algorithm 1: ℓ_1 Tracking**Input:**

- Current frame $F_k \in \mathbb{R}^{h \times w}$.
- Particles $\mathbf{s}_{k-1}^i, i = 1, 2, \dots, N_s$.
- Templates set $A = [T, E] \in \mathbb{R}^{d_0 \times (N_t + 2d_0)}$.
- Templates' weight vector α associated with T .

begin

Generate new particles $\mathbf{s}_k^i, i = 1, 2, \dots, N_s$ within the PF framework;
for $i \leftarrow 1$ **to** N_s **do**
 Obtain observation \mathbf{y}_i corresponding to \mathbf{s}_k^i ;
 Obtain \mathbf{x} via solving (11) with IP-based methods;
 Calculate residual: $r_i = \|\mathbf{y}_i - T \cdot \mathbf{x}_t\|_2$;
end
 $i^* \leftarrow \underset{1 \leq i \leq N_s}{\operatorname{argmin}}(r_i)$;
Get the observed target $\mathbf{y}_k \leftarrow \mathbf{y}_{i^*}$ and its state
 $\mathbf{s}_k \leftarrow \mathbf{s}_k^{i^*}$;
Update templates T and weights α based on \mathbf{x}_{i^*} as in [1];

end**Output:**

- Tracked target \mathbf{y}_k .
- Updated target dynamic state \mathbf{s}_k .
- Updated target templates T and their weights α .

C. Compressed sensing and its application in pattern recognition

CS states that a η -sparse¹ signal $\mathbf{x} \in \mathbb{R}^n$ can be exactly recovered with overwhelming probability via few measurements

$$\mathbf{y}_i = \Phi_i \mathbf{x}, \quad i = 1, \dots, m \ll n.$$

Intuitively, one would achieve \mathbf{x} via

$$\min_{\mathbf{x}} \|\mathbf{x}\|_0, \text{ s.t. } \Phi \mathbf{x} = \mathbf{y}, \quad (12)$$

where $\Phi \in \mathbb{R}^{m \times n}$ is the measurement matrix, of which rows are the measurement vectors Φ_i and $\mathbf{y} = (y_1, \dots, y_m)^T$. $\|\mathbf{x}\|_0$ is the number of non-zero elements of \mathbf{x} . Since (12) is NP-hard [22], it is commonly relaxed to

$$\min_{\mathbf{x}} \|\mathbf{x}\|_1, \text{ s.t. } \Phi \mathbf{x} = \mathbf{y}, \quad (13)$$

which can be casted into a linear programming problem.

As regards CS-based pattern recognition, to deal with noise, one could alternatively solve a Second Order Cone Program:

$$\min_{\mathbf{x}} \|\mathbf{x}\|_1, \text{ s.t. } \|\Phi \mathbf{x} - \mathbf{y}\|_2 \leq \varepsilon, \quad (14)$$

where ε is a pre specified tolerance.

III. REAL-TIME COMPRESSED SENSING TRACKING

In this section, we present the proposed real-time CS tracking.

A. Dimension reduction

The biggest problem of ℓ_1 tracking is the extremely high dimensionality of the feature space, which leads to heavy computation. More precisely, suppose that the cropped image of observation is $I \in \mathbb{R}^{h \times w}$, the dimensionality $d_0 = h \cdot w$ is typically in the order of $10^3 \sim 10^5$, which prevents tracking from real-time.

Fortunately, in the context of compressed sensing (ignoring the non-negativity constraint on \mathbf{x} for now), it is well known that if the measurement matrix Φ has Restricted Isometry Property (RIP) [19], then a sparse signal \mathbf{x} can be recovered from

$$\min \|\mathbf{x}\|_1, \text{ s.t. } \|\Phi A \mathbf{x} - \Phi \mathbf{y}\|_2 \leq \varepsilon. \quad (15)$$

A typical choice of such measurement matrix is random gaussian matrix

$$R \in \mathbb{R}^{d \times n}, \quad R_{i,j} \sim \mathcal{N}(0, 1).$$

Besides random projection, there are other means that guarantee RIP. Shi *et al.* [23] proposed a hash kernel to deal with the issue of computational efficiency. Let $h_s(j, d)$ denotes a hash function (*i.e.*, the hash kernel) $h_s : \mathbb{N} \rightarrow \{1, \dots, d\}$ drawn from a distribution of pairwise independent hash functions, where $s \in \{1, \dots, S\}$ is the seed. Different seed gives different hash function. Given $h_s(j, d)$, the hash matrix H is defined as

$$H_{ij} := \begin{cases} 2h_s(j, 2) - 3, & h_s(j, d) = i, \forall s \in \{1, \dots, S\} \\ 0, & \text{otherwise.} \end{cases} \quad (16)$$

Obviously, $H_{ij} \in \{0, \pm 1\}$. The hash kernel generates hash matrices more efficiently than conventional random matrices while maintains the similar random characteristics, which implies good RIP.

¹a signal \mathbf{x} is said η -sparse if there are at most η nonzero entries in \mathbf{x} .

In this work, the dimensionality of feature space is reduced by matrix $\Phi \in \mathbb{R}^{d \times d_0}$ (which could be either random matrix R or hash matrix H) from d_0 to d where $d \ll d_0$. This significantly speeds up solving equation (14), for its complexity depends on d polynomially.

B. Customized orthogonal matching pursuit for real-time tracking

1) *Orthogonal matching pursuit*: Before the compressed sensing theory was proposed, numerous approaches had been applied for sparse approximation in the literature of signal processing and statistics [24–26]. Orthogonal Matching Pursuit (OMP) is one of the approaches and solves (12) in a greedy fashion. Tropp and Gilbert [22] proved OMP’s recoverability and showed its higher efficiency compared with linear programming which is adopted by the original ℓ_1 tracker of [1]. Be more explicit, given that $A \in \mathbb{R}^{d \times n}$ the computational complexity of linear programming is around $O(d^2 n^{\frac{3}{2}})$, while OMP can achieve as low as $O(dn)^2$. We implement the sparse recovery procedure of the proposed tracker with OMP so as to accelerate the tracking process.

The number of measurements required by OMP is $O(\eta \log(n))$ for η -sparse signals, which is slightly harder to achieve compared with that in ℓ_1 minimization. However, it is merely a theoretical bound for signal recovering, no significant impact of OMP upon the tracking accuracy is observed in our experiments (see Section V).

2) *Further acceleration—OMP with early stop*: The OMP algorithm was proposed for recovering sparse signal exactly (see Equation (12)), and the perfect recovery is also guaranteed within d steps [25]. However, in the realm of pattern recognition, we argue that there is no requirement for perfect recovery for many applications. For example, for classification problems, test accuracy is of interest and exact recovery does not necessarily translate into high classification accuracy. So on the contrary, an appropriate recovery error may even improve the accuracy of recognition [17]. We introduce a residual based stopping criterion into OMP by modifying (12) as

$$\min_{\mathbf{x}} \|\mathbf{x}\|_0, \text{ s.t. } \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2 \leq \varepsilon. \quad (17)$$

Moreover, the procedure of OMP could be accelerated remarkably if the above stopping criterion is enforced. To understand this, let us assume that OMP follows the MP algorithm [24] with respect to the convergence rate³, i.e.,

$$r_k = \frac{K}{\sqrt{t}}, \quad t < n, \quad (18)$$

where K is a positive constant and $r_k = \|\mathbf{A}\mathbf{x}_k - \mathbf{y}\|_2$ is the recovery residual after t steps. Given that we relax the stopping criterion ε by 10 times

$$\varepsilon' = 10\varepsilon, \quad (19)$$

²Here, however, we do not employ the trick that Tropp and Gilbert mentioned for the least-squares routine. As a result, the OMP’s complexity is higher than $O(dn)$ but still much lower than that of linear programming.

³Although the convergence rate for MP algorithm is $O(1/\sqrt{t})$, the convergence rate for OMP remains unclear.

then the required step t_{stop} is reduced to be

$$\begin{aligned} t'_{\text{stop}} &= K^2/\varepsilon'^2 \\ &= 10^{-2}K^2/\varepsilon^2 \\ &= 10^{-2}t_{\text{stop}}. \end{aligned} \quad (20)$$

Considering that the complexity of OMP is at least proportional to t , the algorithm could be accelerated by 100 times theoretically. Figure 1 shows the empirical influence of the terminating criterion upon the running iterations and running time. In our algorithm, we empirically set the stopping threshold $\varepsilon = 0.01$, which draws a balance between speed and accuracy.

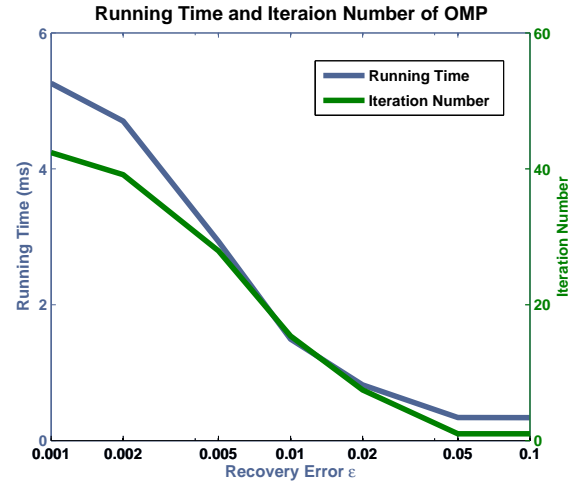


Fig. 1: The decreasing tendency of running time and iteration numbers of the OMP procedure with different residual thresholds. The result is produced from a Matlab-based experiment on video “Cubicle”, with the feature dimension of 50. Both the running time and iteration numbers are the average result over all the frames and particles.

3) *Tracking with a large number of templates*: One noticeable advantage of the SRC-based tracker is the exploitation of multiple templates obtained from different frames. However, for the ℓ_1 tracker, the number of templates n should be curbed into strictly because it equals to the dimensionality of the optimization variable \mathbf{x} . To design a good ℓ_1 tracker, a trade-off between n and the optimization speed is always required. Fortunately, this dilemma dose not exist when the tracker is facilitated with OMP and a carefully-selected sparsity η .

The computational burden of OMP consists of two steps: one is for selecting the maximum correlated vector from matrix $A \in \mathbb{R}^{d \times n}$, and the other is for solving the least squares fitting. In step t ($t < d$), it is trivial to compute the complexity of the first step is $O(dn)$ and that for least-square fitting is $O(d^3 + td^2 + td)$. Accordingly, the running time of OMP is dominated by solving the least-squares problem, which is independent of the number of templates, n . In other words, *within a certain number of iterations, the amount of templates would not affect the overall running time significantly*. This is an important and desirable property in the sense that we might be able to employ a large amount of templates.

Admittedly, larger n might lead to more iterations. However, if we impose a maximum sparsity η , the OMP procedure would only last for η steps in the worst scenario. From this

perspective, a preset $\eta \ll n$ is capable to eliminate the influence of a large n upon the running iterations. Figure 2 depicts the change tendency of running time with increasing n , given that $d \in \{50, 75\}$, $\eta = 15$. As can be seen, the elapsed time is only doubled when n is raised by 10^2 times.

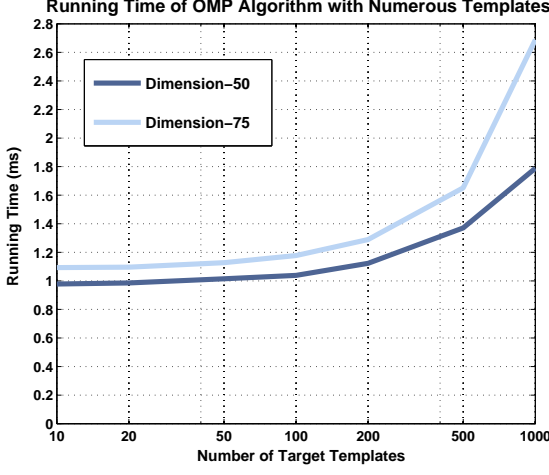


Fig. 2: Running time of OMP with various numbers of target templates. The experiment is carried out on video sequence “Cubicle” with reduced dimensions 50 and 75. The recorded running time is the average time consumption for one OMP procedure which calculates the observation likelihood for a particle. Note that the x -axis only indicates target templates’ number, and the number of trivial templates is not counted. The sparsity $\eta = 15$.

Inspired by this valuable finding, we aggressively set the number of target templates to 100 which is 10 times larger than that in X. Mei’s paper. We try to harness the enormous target templates to accommodate the variation of illumination, gesture and occlusion and consequently improve the tracking accuracy. As regards the sparsity, we elaborately set $\eta = 0.5 \cdot d$ for RTCST and $\eta = 15$ for RTCST-B which is introduced in Section IV. We believe the numbers are sufficiently large for the representations.

Hereby, we sum up all the adjustments to OMP mentioned in Algorithm 2. Note that here we use the inner product rather than its absolute value to verify the correlation. This heuristic manner is used to make the recovered coefficient vector $\mathbf{x} \succeq 0$, approximately. For RTCST-B introduced in next section, the absolute value of inner product is re-employed due to the absence of the positive constraint.

C. Minor modifications

Besides the dimension reduction methods and OMP, modifications to the original ℓ_1 tracker are proposed in this section to achieve a even higher tracking accuracy.

1) *Update templates according to sparsity concentration index:* In the ℓ_1 tracker, the template set is updated when a certain threshold of similarity is reached, i.e.,

$$\text{sim}(\mathbf{y}, \mathbf{a}_i) < \tau, \quad (21)$$

where $i = \arg\max(x_i)$ and $\text{sim}(\mathbf{y}, \mathbf{a})$ is the function for evaluating the similarity between vectors \mathbf{y} and \mathbf{a} . It can be the angle between two vectors or SSD between them. However, Wright *et al.* proposed a better approach to validate the representation. The approach, which utilizes the recovered \mathbf{x} itself

Algorithm 2: Customized OMP for Tracking

Input:

- A normalized observation $\mathbf{y} \in \mathbb{R}^d$.
- A mapped templates set $\Phi A = [\mathbf{a}_1, \dots, \mathbf{a}_n] \in \mathbb{R}^{d \times n}$.
- A recovery residual $0 < \varepsilon \ll 1$.
- A sparsity $0 < \eta \ll n$.

begin

Initialize the residual $\mathbf{r}_0 = \mathbf{y}$, index set $\Lambda_0 = \emptyset$ and selected template set $\Psi_0 = \emptyset$;

for $t \leftarrow 1$ **to** η **do**

$\lambda_t = \arg\max_{j=1, \dots, n} \langle \mathbf{r}_{t-1}, \mathbf{a}_j \rangle$;

$\Lambda_t = \Lambda_{t-1} \cup \{\lambda_t\}$;

$\Psi_t = [\Psi_{t-1} \ \mathbf{a}_{\lambda_t}]$;

Solve the least-squares problem:

$$\mathbf{x}_t = \arg\min_{\mathbf{x}} \|\Psi_t \mathbf{x} - \mathbf{y}\|_2;$$

Calculate the new residual:

$$\mathbf{r}_t = \mathbf{y} - \Psi_t \mathbf{x}_t;$$

if $\|\mathbf{r}_t\|_2 < \varepsilon$ **then break**;

end

Retrieve signal \mathbf{x} according to \mathbf{x}_t and Λ_t ;

end

Output:

- Recovered coefficients $\mathbf{x} \in \mathbb{R}^n$
-

rather than the similarity, is termed *Sparsity Concentration Index* (SCI) [17]. Particularly, in the context of RTCST, class number is 1 if the noise is not viewed as a class, then we obtain a simplified SCI measurement for the target class, which writes

$$\text{SCI}_t(\mathbf{x}) = \|\mathbf{x}_t\|_1 / \|\mathbf{x}\|_1 \in [0, 1], \quad (22)$$

where $\mathbf{x}_t = \mathbf{x}(1 : N_t)$. In the presented RTCST algorithm, SCI_t is employed instead of (21).

2) *Abandoning the template weight:* The original ℓ_1 tracker enforces a template re-weighting scheme to distinguish templates by [1], their importance. Nonetheless, following their scheme the weight of each target template is always smaller than that of noise templates (see Algorithm 1). This does not make much sense. Actually, it may be intractable to design an ideal template re-weighting scheme that works in all the circumstances. A poorly-designed re-weighting scheme could even deteriorate the tracking performance. We abandon the template weight because the importance of templates be easily exploited by the compressed sensing procedure. Without template weights, the tracker becomes simpler and less heuristic. The empirical result also shows better tracking accuracy when template weight is abandoned.

3) *MAP and MSE:* In Mei and Ling’s framework [1], the new state \mathbf{s}_k is corresponding to the particle with the largest observation likelihood. This method is known as the Maximum A Posterior (MAP) estimation. It is also known that for the particle filtering framework, Mean Square Error (MSE) estimation is usually more stable than MAP. As a result, we adopt MSE in our real-time tracker, namely,

$$\mathbf{s}_k = \frac{\sum_{i=1}^{N_s} (\mathbf{s}_k^i \cdot l_i)}{\sum_{i=1}^{N_s} l_i}, \quad (23)$$

Algorithm 3: Real-Time Compressed Sensing Tracking**Input:**

- Current frame $F_k \in \mathbb{R}^{h \times w}$.
- Particles $\mathbf{s}_{k-1}^i, i = 1, 2, \dots, N_s$
- A dimension-reduction matrix $\Phi \in \mathbb{R}^{d \times d_0}$
- A Templates set $A = [T, E] \in \mathbb{R}^{d_0 \times (N_t + 2d)}$.
- A preset parameter $\lambda > 0$.

beginNormalize every column of ΦA ;Generate new particles $\mathbf{s}_k^i, i = 1, 2, \dots, N_s$;**for** $i \leftarrow 1$ **to** N_s **do**Obtain mapped observation $\Phi \mathbf{y}_i$ corresponding to \mathbf{s}_k^i ;Get \mathbf{x} via solving (24) with Algorithm 2;Calculate residual r_i via (25);Calculate observation likelihood $l_i = \exp(-\lambda \cdot r_i)$ **end**Calculate target dynamic state \mathbf{s}_k via (23) and then get the target \mathbf{y}_k ;Recalculate \mathbf{x}_k for \mathbf{y}_k via solving (24);Update templates T based on \mathbf{x}_k and (22);**end****Output:**

- Tracked target \mathbf{y}_k .
- Updated target dynamic state \mathbf{s}_k .
- Updated target templates T .

where \mathbf{s}_k^i is the i th particle at time k and l_i is the corresponding observation likelihood.

D. The Algorithm

In a nutshell, for each observation, we utilize Algorithm 2 to recover the coefficient vector \mathbf{x} by solving the problem

$$\min_{\mathbf{x}} \|\mathbf{x}\|_0, \text{ s.t. } \|\Phi A \mathbf{x} - \Phi \mathbf{y}\|_2 \leq \varepsilon, \mathbf{x} \succeq 0 \quad (24)$$

where $\mathbf{x} = [\mathbf{x}_t, \mathbf{x}_e]$, $A = [T, E]$. The residual is then obtained by

$$r = \|\Phi \mathbf{y} - \Phi A \mathbf{x}_t\|_2. \quad (25)$$

Finally the likelihood of this observation is updated as

$$l = \exp(-\lambda \cdot r), \lambda > 0. \quad (26)$$

The procedure of Real-Time Compressed Sensing Tracking algorithm is summarized in Algorithm 3. Our template update scheme is demonstrated in Algorithm 4. As can be seen, the proposed update scheme is much conciser than that in the ℓ_1 tracker [1] thanks to the abandonment of template weight. The empirical performance of RTCST is verified in Section V.

IV. RTCST-B: MORE ROBUST AND EFFICIENT RTCST WITH BACKGROUND MODEL

To some extent, visual tracking is viewed as object detection task with prior information. Similar to object detection, which is sometimes treated as a classification problem, visual tracking also distinguishes the foreground (target) from background. In detection applications, the background class

Algorithm 4: Template Update Scheme for RTCST**Input:**

- Sparse coefficient $\mathbf{x} = \mathbf{x}_k$ in Alg. 3.
- Observed target \mathbf{y}_k .
- Target templates set $A = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_{N_t}] \in \mathbb{R}^{d_0 \times N_t}$.
- A preset parameter $0 < \tau < 1$.

begin**if** $SCI_t(\mathbf{x}) < \tau$ **then** $j^* \leftarrow \underset{1 \leq j \leq N_t}{\operatorname{argmin}}(\mathbf{x}_j);$ $\mathbf{a}_{j^*} \leftarrow \mathbf{y}_k$, where \mathbf{a}_{j^*} is the j^* -th target template;**end****end****Output:**

- Updated target templates A .

is usually considered without distinct feature because it could follow any pattern. Quite the contrary, in the context of visual tracking, the background is much more limited with respect to appearance variation. Particularly, for the stationary camera, the background is nearly fixed. Under these assumptions, it is worthwhile exploiting the background information for tracking. And appropriate incorporations of background model indeed improve the tracking performance[7, 27–29].

We hereby propose a novel CS-based background model (CSBM) to facilitate tracking algorithm. The definition of CS-based background model is quite simple. Suppose that $\Gamma_i \in \mathbb{R}^{h \times w}$, $i = 1, \dots, N_b$ is the i th frame where foreground is absent, and h and w are the height and width of the frame respectively, we define the background model as

$$\mathbb{G} = \{\Gamma_1, \dots, \Gamma_{N_b}\} \quad (27)$$

or in short, the collection of N_b backgrounds. The background templates are then generated from CSBM to cooperate with target templates in our new tracker.

Please note that our algorithms is unrelated to the background subtraction manner proposed by Volkan *et al.* [20]. In their paper, foreground silhouettes are recovered via CS procedure but the background subtraction is still performed in conventional way. Our CSBM and RTCST-B is entirely different from their manner, both in essence and appearance. The details of CSBM and its incorporation with RTCST are introduced below.

A. Building the Optimal CSBM

A good CSBM should only constitute “pure” backgrounds and contain sufficiently large appearance variation, *e.g.*, illumination changes. Ideally, we could simply select certain number of foreground-absent frames from video sequence to build a CSBM. However, the “pure” background is usually difficult to find and it is even harder to obtain the ones cover the main distribution of background appearance.

An intuitive way to obtain a clean background is replacing the foreground of one frame with a background patch cropped from another frame. More precisely, let $F \in \mathbb{R}^{h \times w}$ denote the frame based on which the background is retrieved, and

$F' \in \mathbb{R}^{h \times w}$ stand for the frame where the background patch is cropped, suppose that the foreground region in F is $F(t : b, l : r)^4$, the patching operation could be described as

$$\Gamma_{i,j} = \begin{cases} F'_{i,j}, & t \leq i \leq b \ \& \ l \leq j \leq r \\ F_{i,j}, & \text{otherwise.} \end{cases} \quad (28)$$

where Γ is the retrieved background. An illustration of (28) is also available in Figure 3.

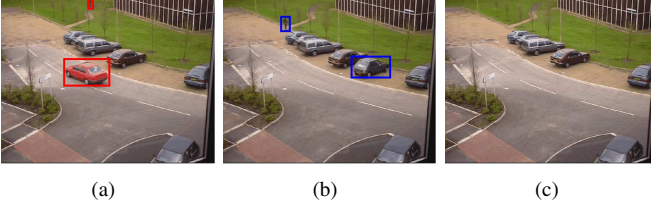


Fig. 3: An illustration for retrieving background. (a) shows contaminated background with foreground regions signed by red rectangles; (b) is the frame where the background patches are obtained, note that blue rectangles indicate the foregrounds in (b), they are far from the foreground areas in (a); (c) demonstrates a retrieved background based on (a) and (b). The frames are captured from video sequence *pets2000_c1*.

In practice, multiple foreground regions need to be mended for each “impure” background candidate. Furthermore, a selection approach should be conducted to form the optimal combination of the retrieved backgrounds over all the candidates. To achieve this goal, we first randomly capture $N' > N_b$ frames from the concerned video sequence. Afterwards, every foreground region of the frames are located manually. The foreground is then replaced by a clean background region cropped from the nearest frame (in terms of frame index). Finally, a k -median clustering algorithm is carried out for selecting N_b most comprehensive backgrounds.

It is nontrivial to notice that even some backgrounds are not perfectly retrieved, *i.e.*, with minor foreground remains, CSBM can still work well considering that CS is robust to the noise in measurements [19].

B. Equipping RTCST with CSBM

We equip the RTCST with CSBM to build a novel visual tracker, *a.k.a Real-Time CS-based Tracker with Background Model* (RTCST-B). In RTCST-B, original noise templates are replaced by *background templates* which are generated from CSBM. In the context of PF tracking, given a observation position Ξ with d_0 pixels and a CSBM \mathbb{G} defined in (27), the background templates set B is obtained by:

$$B = [I_1 \ I_2 \ \dots \ I_{N_b}] \in \mathbb{R}^{d_0 \times N_b} \quad (29)$$

$$I_i = \text{CV}(\Gamma_i, \Xi) \ \forall i = 1, \dots, N_b$$

where function $\text{CV}(\cdot)$ is called *crop-vectorize* operation which first crops the region indicated by Ξ from background Γ_i and then vectorize it into $I_i \in \mathbb{R}^{d_0}$. Eventually, the optimization problem for RTCST-B writes:

$$\min_{\mathbf{x}} \|\mathbf{x}\|_0 \ \text{s.t.} \ \|\Phi \mathbf{A} \mathbf{x} - \Phi \mathbf{y}\|_2 \leq \varepsilon, \quad (30)$$

⁴In this paper, all the target or foreground is represented as a rectangle region

where \mathbf{x} is comprised of \mathbf{x}_t and \mathbf{x}_b , *i.e.*, the coefficient vectors for target and background, $A = [T, B] \in \mathbb{R}^{d \times (N_t + N_b)}$.

Despite the diverse optimization problem, the calculation for the likelihood remains the same as in (26). To understand this, let \mathbf{x}_t and \mathbf{x}_b denote the coefficients associated with target templates and background templates respectively, $p(\mathbf{y}_k | \mathbf{s}) = p(\mathbf{y}_k | \mathbf{x}_t) = \exp(-\lambda r)$ be the observation likelihood⁵, where r is defined in (25), then we have:

$$p(\mathbf{y}_k | \mathbf{x}_t, \mathbf{x}_b) = p(\mathbf{y}_k | \mathbf{x}_t) = \exp(-\lambda r) \quad (31)$$

with the assumption that \mathbf{x}_t and \mathbf{x}_b are deterministic by each other, *i.e.*,

$$p(\mathbf{x}_b, \mathbf{x}_t) = p(\mathbf{x}_t) = p(\mathbf{x}_b) \quad (32)$$

or in other words, the solution of CS procedure is unique. [19].

In addition, the template update scheme should be changed slightly considering a new class is involved in. More precisely, target templates are updated only when

$$\text{SCI}_{tb}(\mathbf{x}) = \frac{\max\{\|\mathbf{x}_t\|_1, \|\mathbf{x}_b\|_1\}}{\|\mathbf{x}\|_1} \leq \tau \quad (33)$$

Finally, the positive constraint for \mathbf{x} is removed in 30 because background subtraction implies minus coefficients for background templates. It is reasonable to not curb the coefficients in RTCST-B.

In summary, one just needs to impose following minor modifications on RTCST to transfer it into RTCST-B.

- 1) Substitute the background templates for noise templates.
- 2) Eliminate the positive constraint.
- 3) Conduct the CV operation for each observation.
- 4) Utilize the new SCI measurement.

Apparently, the diversity between RTCST and RTCST-B is not significant with respect to formulation. Nevertheless, the seemingly small change makes RTCST-B much more superior to its prototypes.

C. Superiority Analysis

Compared with the ℓ_1 tracker and RTCST, RTCST-B enjoys three main advantages which are described as follows.

1) *More Sparse*: An underlying assumption behind the ℓ_1 tracker and RTCST is that, the background could be sparsely represented by noise templates in E . It is true when foreground dominates the observed rectangle. More quantitatively, given η_t is the sparsity of target coefficient vector \mathbf{x}_t , when

$$\eta_t + \|\mathbf{x}_e\|_0 \leq d/3$$

the representation based on solution \mathbf{x} in (24) is guaranteed to be reliable [17]. Nonetheless, the sparse representation is no longer valid when the background covers the main part of observation. Predictably, the incorrect representation will deteriorate tracking accuracy.

On the other hand, after noise templates being replaced by background templates, the aforementioned assumption usually keeps true. Figure 4 give us a explicit demonstration for the sparsity of solutions.

⁵It is trivial to prove that the relationship between particle \mathbf{s} and \mathbf{x}_t is deterministic given a specific frame image

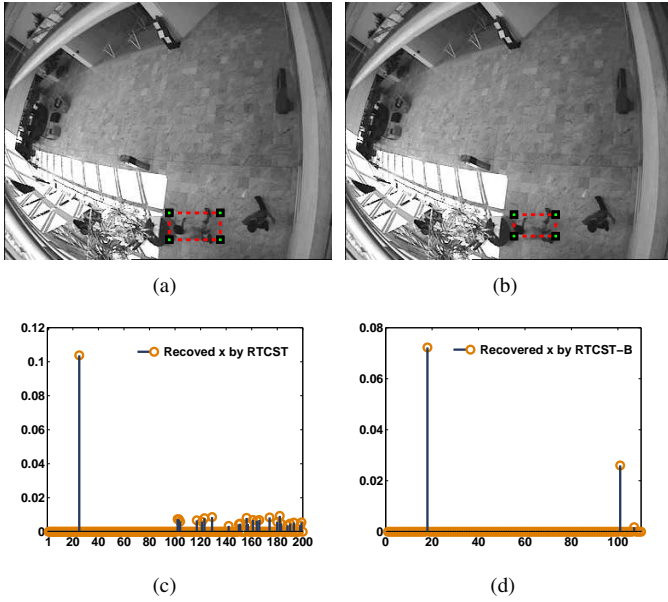


Fig. 4: A demonstration of the sparse solutions for RTCST and RTCST-B. (a) and (b) are the tracking result by RTCST and RTCST-B on the same frame (captured from *pets2004_p1*). (c) and (d) are the recovered signals for RTCST and RTCST-B respectively. The representation by RTCST-B is much more sparse than that by RTCST. Note that here, $d = 50$, $N_t = 100$ and $N_b = 10$ for RTCST-B.

2) *More Efficiency*: Comparing with existing background models, the computation burden of CSBM is extremely trivial. First of all, there is no need to conduct the background subtraction or foreground connection in RTCST-B, because these two functions are integrated within the CS procedure implicitly. Secondly, if the CSBM is generated properly, *i.e.*, can cover the main distribution of background's appearance, to update model becomes unnecessary. Thirdly, the sufficient number of background templates is much smaller than that of noise template, *i.e.*,

$$N_b \ll N_n = 2d$$

where N_n is the number of noise templates. The reduction of templates' amount will immediately speed up the optimization process. The last, and the most important reason is, the required sparsity η for RTCST-B is much smaller than that for RTCST (see Section III-B3). This leads to an earlier terminated OMP procedure in RTCST-B and hence makes it faster. In conclusion, the introduction of CSBM won't impose further computational burden on the algorithm, and just the opposite, the tracking procedure will be accelerated to some extent.

3) *More Robust*: In RTCST and ℓ_1 tracker, one tries to use noise templates $E = [I - I]$ to represent background. However, it is the columns in I , which is called *standard basis vectors*, doesn't favor background images over targets. This character makes RTCST and ℓ_1 tracker powerless for recognizing background and consequently, decreases the tracking accuracy. Differing from the prototype, RTCST-B harnesses the discriminant nature of CS-based pattern recognition. Both foreground (target) and background are treated as a typical class with distinct features. In RTCST-B, target templates compete against background templates, who are as powerful as their competitors, to "attract" the observation. Intuitively,

the more discriminative templates will make RTCST-B more robust.

Moreover, once the tracked region drifts away, background information would be brought into target templates via template update (which is almost unavoidable). In this situation, for RTCST and ℓ_1 tracker, some target templates could be more similar to background than all the noise templates. This leads to a serious classification ambiguity and therefore, poor tracking performance. Quite the contrary, RTCST-B could draw back the target to the correct position thanks to the capacity of recognizing background. In plain words, RTCST-B always tends to locate the target in the region which doesn't *look like* background. An empirical evidence for the robustness of RTCST-B is shown in Figure 5.

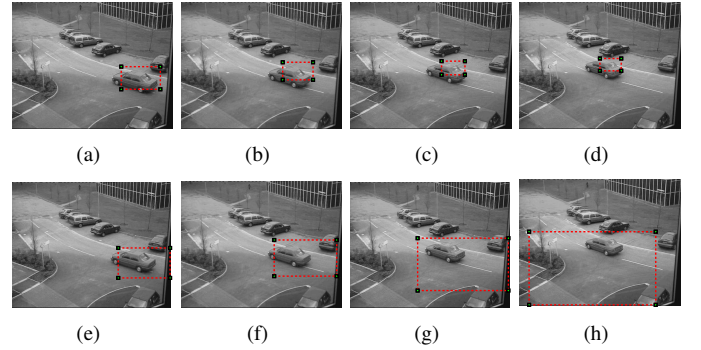


Fig. 5: An empirical evidence for the robustness of RTCST-B against drift. (a) to (d) are the tracking results for RTCST-B compared with image (e) to (h) which are the results for RTCST on the same frames from *pets2000*. The tracked target is signed by red rectangle. We can see that a drift tendency shown on (b) is curbed in the successive frames. Quite the contrary, in the bottom line, the drift effect grows dramatically.

V. EXPERIMENT

A. Experiment Setting

To verify the proposed tracking algorithms, we design a series of experiments for examining the tracking algorithm in terms of accuracy, efficiency and robustness. The proffered algorithms are conducted on 10 video sequences comparing with ℓ_1 tracker, Kernel-Mean-Shift (KMS) tracker [14] and color-based PF tracker [13]. The details of selected video sequences are list in Table II. Note that we only conduct ℓ_1 tracker on 5 videos which are *cubicle*, *dp*, *car11*, *pets2001_c1* and *pets2004-2_p1* respectively. It is because for other videos, the convex optimization problem is too slow to be solved (above 5 minutes per frame).

There are two alternative dimension-reduction manners for RTCST and RTCST-B, namely, random projection and hash matrix projection. In our experiments, both of them are performed with reduced dimension 25, 50 and 100. As regards the particles' number, we examine the proposed trackers with 100 and 200 particles and the numbers for PF tracker is 100, 200 and 500. All the PF-based trackers are run for 20 times except ℓ_1 tracker which is merely conducted for 3 times. We perform KMS tracker for only 1 time considering it is a deterministic method. The average values and standard errors are reported in this section. The MS tracker, PF tracker and ℓ_1 tracker are implemented in C++ while our CS-based trackers

TABLE II: The details of video sequences which are employed for our experiment. The tracking frames refer to the concerned frame index for each video; initial position indicates the minimum bounding box for the target in the first frame; if “Yes” shows in the last column, the video is captured from a stationary camera and consequently, it suits RTCST-B.

	tracking frames	initial position	stationary camera
cubicle	1 ~ 51	[56, 24, 90, 67]	No
dp	1 ~ 66	[91, 25, 116, 57]	No
car4	1 ~ 300	[139, 102, 356, 283]	No
car11	1 ~ 393	[69, 123, 104, 157]	No
fish	1 ~ 200	[122, 57, 208, 148]	No
pets2000_c1	122 ~ 312	[536, 318, 743, 432]	Yes
pets2001_c1	1550 ~ 1635	[8, 272, 46, 296]	Yes
pets2002_p1	275 ~ 500	[578, 92, 641, 172]	Yes
pets2004_p1	115 ~ 550	[193, 258, 251, 287]	Yes
pets2004-2_p1	1 ~ 201	[181, 224, 239, 262]	Yes

are implemented in Matlab. To compare the efficiency with the proposed algorithms, there is also a Matlab version of ℓ_1 tracker. All the algorithms are run on a PC with 2.6GHz quad-core CPU and 4G memory (we only use one core of it). As to the software, we use Matlab 2009a and the linear programming solver is called from Mosek 6.0[30].

It is important to emphasize that in our experiment, *no trick is used for selecting the target region in the first frame*. The initial target region is always the minimum rectangle $R = [l, r, t, b]$ which can cover the whole target⁶, where l , r , t , and b are the left, right, top and bottom boundaries’ coordinates (horizontal or vertical) respectively. This rigid rule is followed for eliminating the artificial factors in visual tracking and making the comparison unprejudiced.

B. TSP — A New Metric of Tracking Robustness

A conventional choice of the manner to verify the tracking accuracy is *tracking error*. Specifically, given that the centroid of ground truth region is \mathbf{c}_g while that of tracked region is \mathbf{c}_t , the tracking error ρ is defined as

$$\rho = \|\mathbf{c}_g - \mathbf{c}_t\|_2, \quad (34)$$

i.e., the euclidean distance between two centroids. However, if we take scale variation into consideration, ρ is poor to verify tracker’s performance. Let’s see Figure 6(a) for an example. In the image, red rectangle indicates the ground truth for a moving car. The blue and gray rectangles, which are obtained by various tracking algorithms, share the identical centroid. By using tracking error, same performance is reported for both two trackers despite the obvious difference on tracking accuracy.

Inspired by the evaluation manner proposed for PASCAL data base[31], we propose a new tracking accuracy measurement which is termed *Tracking Success Probability* (TSP). To obtain the definition of TSP, firstly let’s suppose the bounding box of ground truth region is $R_g = [l_g, r_g, t_g, b_g]$, and the one for tracked region is $R_t = [l_t, r_t, t_t, b_t]$. We then design a function $a(R_g, R_t) \in [-1, 1]$ to estimate the overlapping state between R_g and R_t . Given two distance sets:

$$\begin{aligned} \mathbb{H} &= \{r_t - l_g, r_g - l_t, r_g - l_g, r_t - l_t\} \\ \mathbb{V} &= \{b_t - t_g, b_g - t_t, b_g - t_g, b_t - t_t\} \end{aligned}$$

⁶Shadows are not taken into consideration.

and an indicator function s_{tg}

$$s_{tg} := \begin{cases} -1, & R_g \text{ and } R_t \text{ are separate} \\ 1, & \text{otherwise.} \end{cases} \quad (35)$$

then $a(R_g, R_t)$ writes⁷

$$a(R_g, R_t) = s_{tg} \cdot \left| \frac{\min(\mathbb{H}) \cdot \min(\mathbb{V})}{\max(\mathbb{H}) \cdot \max(\mathbb{V})} \right|,$$

It is easy to find that when two regions overlap each other, $a(R_g, R_t)$ is the ratio of the intersection area $R_{g \cap d}$ to the area R^* , which is the minimum region covering both R_g and R_d . See Figure 6(b) for an instance. Finally, TSP is formulated as

$$\text{TSP}(R_g, R_t) = \frac{\exp(\nu \cdot a(R_g, R_t))}{1 + \exp(\nu \cdot a(R_g, R_t))} \in [0, 1], \quad (36)$$

where $\nu > 0$ is a preset parameter reflects the worst scenario we could assure the target is located correctly. In our experiment, ν is the solution of

$$\frac{\exp(0.25\nu)}{1 + \exp(0.25\nu)} = 0.95 \implies \nu = 11.8. \quad (37)$$

In other words, when the overlapped region is larger than 25% part of region R^* , we are convinced (with the probability of 0.95) that the tracking is successful.

Obviously, the larger the TSP is, the more confident we believe this tracking is successful. If we apply TSP to the tracking results shown in Figure 6(a), then the TSP of blue rectangle is 0.95 which is significantly larger than that of the gray one (with TSP of 0.55). The difference implies that TSP is capable to accommodate dynamic factors besides displacement. Another merit of TSP is the comparability over different video sequences thanks to its fixed value range i.e., $[0, 1]$. Considering these advantages, in the current paper, all the empirical results are evaluated by TSP. As a reference, tracking error results are also available.

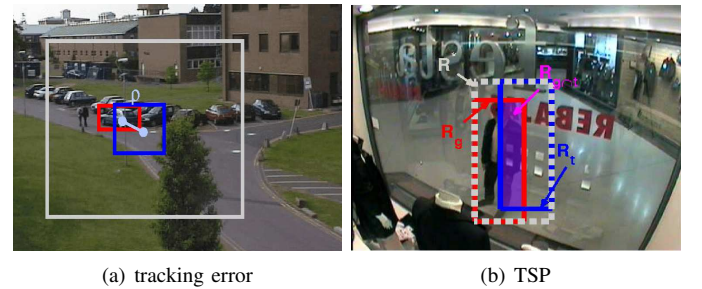


Fig. 6: A demonstration of two measurements of tracking accuracy. (a) shows the poor capacity of ρ . (b) illustrates the definition of TSP. R_g and R_t are illustrated as red and blue rectangles respectively; the region R^* is a gray dashed square in the image while intersection region $R_{g \cap t}$ is shown in purple. We can see that in this case, $a(R_g, R_d) = R_{g \cap t} / R^*$. These two frames are obtained from video sequence *pets2001* and *pets2002* respectively.

C. Tracking Accuracy

Firstly, we examine the tracking accuracy of our trackers comparing with the competitors. The average TSP for every experiment is shown in Table III. For each video sequence, the optimal accuracy is displayed in bold type.

⁷Here, we suppose the origin of image is on the left-top corner.

TABLE III: TSP values for tracking experiments. The term “R-D x -Rand” stands for RTCST with x -dimension features which is generated by random projection while the row started with “RB-...” refers to the results with RTCSTB. “PN x ” indicates x particles are used in the tracker. The optimal values for each video sequence is illustrated in bold type.

		cubicle	dp	car4	car11	fish	pets2000_c1	pets2001_c1	pets2002_p1	pets2004_p1	pets2004-2_p1
KMS		77 \pm 32.7	100 \pm 0.4	24 \pm 29.5	67 \pm 40.5	98 \pm 1.4	94 \pm 4.8	23 \pm 14.6	23 \pm 37.0	52 \pm 30.2	26 \pm 32.8
PF	PN100	95 \pm 8.1	98 \pm 3.2	64 \pm 30.0	37 \pm 30.6	90 \pm 14.2	45 \pm 31.5	97 \pm 2.7	24 \pm 38.2	23 \pm 25.8	58 \pm 16.9
	PN200	95 \pm 8.3	98 \pm 3.0	65 \pm 30.3	39 \pm 31.8	90 \pm 14.7	44 \pm 31.7	98 \pm 2.5	24 \pm 38.4	23 \pm 25.9	58 \pm 16.9
	PN500	95 \pm 7.9	98 \pm 2.9	64 \pm 33.6	39 \pm 32.7	90 \pm 15.4	44 \pm 33.1	98 \pm 2.5	24 \pm 38.4	22 \pm 25.8	58 \pm 17.0
R-D25-Rand	PN100	69 \pm 21.4	66 \pm 20.1	89 \pm 11.0	64 \pm 17.2	63 \pm 20.1	77 \pm 8.9	89 \pm 8.3	54 \pm 21.6	31 \pm 25.8	33 \pm 29.1
	PN200	80 \pm 15.8	78 \pm 15.7	95 \pm 7.5	62 \pm 20.7	64 \pm 20.2	80 \pm 8.5	87 \pm 10.1	63 \pm 16.0	28 \pm 25.6	29 \pm 31.1
R-D50-Rand	PN100	73 \pm 21.5	78 \pm 16.7	95 \pm 8.1	64 \pm 24.1	61 \pm 21.0	72 \pm 10.1	86 \pm 12.5	65 \pm 16.1	28 \pm 25.3	25 \pm 33.1
	PN200	69 \pm 23.0	82 \pm 17.9	95 \pm 10.7	81 \pm 22.3	64 \pm 19.0	81 \pm 9.1	83 \pm 13.4	64 \pm 15.1	31 \pm 25.2	25 \pm 32.9
R-D100-Rand	PN100	70 \pm 24.7	71 \pm 21.5	94 \pm 11.2	85 \pm 24.6	64 \pm 19.3	72 \pm 12.5	93 \pm 5.1	61 \pm 16.1	28 \pm 27.2	26 \pm 32.0
	PN200	72 \pm 22.3	77 \pm 17.6	96 \pm 8.8	78 \pm 23.1	59 \pm 20.6	81 \pm 8.7	91 \pm 6.9	68 \pm 13.6	32 \pm 25.7	24 \pm 33.2
R-D25-Hash	PN100	73 \pm 21.3	76 \pm 12.0	90 \pm 12.1	65 \pm 24.4	64 \pm 19.9	83 \pm 6.4	77 \pm 20.3	67 \pm 15.0	38 \pm 25.0	32 \pm 29.6
	PN200	77 \pm 18.3	81 \pm 14.6	89 \pm 14.8	59 \pm 23.9	63 \pm 20.2	96 \pm 2.7	70 \pm 23.5	55 \pm 19.6	35 \pm 25.8	33 \pm 29.8
R-D50-Hash	PN100	73 \pm 21.8	79 \pm 16.5	98 \pm 3.2	75 \pm 24.1	66 \pm 21.3	73 \pm 11.8	100 \pm 0.1	64 \pm 16.0	34 \pm 25.3	22 \pm 32.3
	PN200	75 \pm 21.7	83 \pm 14.2	99 \pm 1.2	74 \pm 22.5	68 \pm 21.6	79 \pm 10.5	100 \pm 0.1	63 \pm 15.7	39 \pm 24.4	21 \pm 33.5
R-D100-Hash	PN100	82 \pm 15.0	88 \pm 10.1	95 \pm 9.1	80 \pm 32.9	56 \pm 22.0	91 \pm 4.6	100 \pm 0.1	64 \pm 14.5	30 \pm 25.9	27 \pm 31.6
	PN200	90 \pm 8.3	92 \pm 8.5	95 \pm 9.3	80 \pm 33.6	52 \pm 23.8	92 \pm 5.3	100 \pm 0.1	67 \pm 13.3	30 \pm 26.3	28 \pm 31.8
RB-D25-Rand	PN100	—	—	—	—	—	76 \pm 7.0	86 \pm 9.5	80 \pm 8.8	68 \pm 18.2	49 \pm 23.2
	PN200	—	—	—	—	—	92 \pm 3.4	84 \pm 12.1	78 \pm 10.2	62 \pm 17.4	59 \pm 19.2
RB-D50-Rand	PN100	—	—	—	—	—	86 \pm 5.8	98 \pm 2.0	73 \pm 11.8	58 \pm 18.4	44 \pm 26.5
	PN200	—	—	—	—	—	93 \pm 3.6	97 \pm 2.7	77 \pm 10.8	58 \pm 18.3	62 \pm 17.8
RB-D100-Rand	PN100	—	—	—	—	—	96 \pm 4.2	100 \pm 0.6	74 \pm 11.6	46 \pm 24.0	54 \pm 20.8
	PN200	—	—	—	—	—	95 \pm 5.0	100 \pm 0.1	72 \pm 11.8	51 \pm 21.7	53 \pm 22.0
RB-D25-Hash	PN100	—	—	—	—	—	89 \pm 2.9	94 \pm 6.1	79 \pm 10.3	64 \pm 20.7	71 \pm 14.4
	PN200	—	—	—	—	—	89 \pm 3.6	89 \pm 8.9	77 \pm 10.5	61 \pm 16.1	77 \pm 12.0
RB-D50-Hash	PN100	—	—	—	—	—	75 \pm 12.0	98 \pm 1.7	82 \pm 9.0	42 \pm 25.3	52 \pm 22.7
	PN200	—	—	—	—	—	98 \pm 1.9	98 \pm 1.7	82 \pm 8.7	59 \pm 19.5	71 \pm 14.0
RB-D100-Hash	PN100	—	—	—	—	—	97 \pm 1.9	99 \pm 1.3	82 \pm 8.9	51 \pm 20.8	67 \pm 14.7
	PN200	—	—	—	—	—	99 \pm 1.4	98 \pm 1.7	82 \pm 9.8	53 \pm 22.2	71 \pm 13.1
LIT		99 \pm 2.2	92 \pm 8.8	—	77 \pm 37.4	—	—	100 \pm 0.0	—	34 \pm 26.4	—

As illustrated in Table III, all the tracking approaches achieve similar performances on the sequence with simple background and stable illumination (*dp* and *cubicle*). For the video sequence *fish*, traditional methods show higher capacity for accommodating extreme illumination variation. On the other hand, for the outdoor scene and complex background tasks, *i.e.*, the other 7 sequences, CS-based trackers consistently outperform PF tracker and KMS tracker. All the best performances are observed with RTCST and RTCST-B for these video sequences. Considering that the target could be viewed as missed when the TSP is below 30%, the traditional trackers are failure for the majority of these video datasets, *i.e.*, KMS tracker for *car4*, *pets2001_c1*, *pets2002_p1* and *pets2004-2_p1*; PF tracker for *pets2002_p1* and *pets2004_p1*. Moreover, ℓ_1 tracker also fails on *pets2004_p1* and *pets2004-2_p1* due to the unstable target appearances. Our methods, on the contrary, do much better than the competitors and handle some intractable sequences (*e.g.*, *pets2004_p1* and *pets2004-2_p1*) very smoothly (with the TSP > 65%). Particularly, for the camera-fixed scenes, RTCST-B is applied and always achieves the highest accuracy. The superiority of RTCST-B over all the other trackers confirms our assumption that higher accuracy would be achieved when the tracking is considered as binary classification problem.

Besides the TSP values, video frames with the tracked regions are listed in Figure 8 while tracking errors changing along with the frame index are also plotted in Figure 9.

In Figure 8, only the best (with the highest average TSP value) result is employed to be shown for each tracker. The explicit tracking results support the statistics in Table III. RTCST beats KMS tracker and PF tracker on *cubicle*, *car4*, *pets2000_c1* and *pets2002_p1* and obtain the similar performance as its competitors on *dp*. Being facilitated with CSBM, RTCST-B always achieves the highest accuracy if it is present. Quite the contrary, the traditional trackers fail in some complex

scenarios, *e.g.* PF tracker on *car4* and *pets2002_p1*; KMS tracker on *car4* and *pets2002_p1*.

From the error curves shown in Figure 9, we can find that our methods beat other visual tracking algorithms on most video sequences except *dp* and *fish*. Given that all the trackers perform similarly for *dp* and video *fish* is generated with extreme illumination variation which is added deliberately, RTCST and RTCST-B could be considered better than their competitors in terms of accuracy.

To evaluate the new measurement, the TSP curves for *cubicle* and *pets2002_p1* are also available in the Figure 9(k) and Figure 9(l). We can see that the TSP value and tracking error change oppositely, which is as expected. However, based on TSP, we can verify the capacity of single tracker without any “reference tracker”. This is hard to achieve based on tracking error.

D. Tracking Efficiency

Efficiency plays a fatal role in real-time visual tracking applications. We record the elapsed time of each tracker in our experiment. The time consumptions (in *ms*) for processing one frame by the tracking algorithms are reported in Table IV. In the table, huge differences in tracking speed are observed. KMS tracker illustrates the highest efficiency with the lowest running speed of 83 *ms per frame* (83 *mspf*). On the contrary, ℓ_1 tracker (both for C-based version and Matlab-based version) is consistently slower than 14000 *mspf* due to the high computational complexity. Being equipped with OMP and dimension reduction manners, RTCST and RTCST-B are able to accelerate the original CS-based tracker by 117.3 (*dp*) to 6271.2 (*pets2004_p1*) times. The speed range for RTCST is 54 \sim 968 *mspf* while that for RTCST-B is 85 \sim 534 *mspf*. PF tracker shows unstable efficiency among all the tests. Its running speed varies from 37 to 1727 *mspf* for the experiment

with 500 particles. Supposed that the speed threshold for real-time application is 100 *mspf*, most of the traditional methods and a part of our methods are qualified. ℓ_1 tracker could not be viewed as “real-time” from any perspective.

Moreover, since RTCST and RTCST-B are implemented in Matlab with single core, their running speeds could be increased remarkably by employing C/C++ language and multiple cores. Actually, the speed of Matlab-based ℓ_1 is already raised by 3.7 (*pets2004_p1*) to 8.4 (*cubicle*) times in its C/C++ counterpart even though only one core is used. If we conservatively predict 10-time speed growth, both RTCST and RTCST-B will be qualified for real-time application in all the circumstances.

E. Tracking Robustness

As mentioned before, no trick is played to select the initial target region. The first region R should always be the minimum bounding box covers the whole target. Nonetheless, the bounding box could merely obtained manually, and hence, approximately. In practice, the selection error is unavoidable. If the visual tracker is not robust enough, minor selection error would lead to massive deviation with respect to tracking performance. We design a new experiment to test the robustness of tracking algorithms. In every repetition of the experiment, a fluctuation vector $\delta = [\delta_l, \delta_r, \delta_s]$, is generated randomly as

$$\delta_l \sim \mathcal{N}(0, \omega), \delta_r \sim \mathcal{N}(0, \omega), \delta_s \sim \mathcal{N}(0, \frac{\omega}{25})$$

where ω is a preset standard deviation with small value. The original bounding box $R = [l, r, t, b]$ is then imposed by δ to obtain a fluctuated rectangle region R^* as

$$R^* = [l^*, r^*, t^*, b^*]$$

where l^* , r^* , t^* and b^* are the new coordinates which are defined as

$$\begin{aligned} l^* &= l + \delta_l, \quad t^* = t + \delta_t, \\ r^* &= (1 + \delta_s) \cdot (r - l) + l + \delta_l, \\ b^* &= (1 + \delta_s) \cdot (b - t) + t + \delta_t. \end{aligned}$$

The tracking is then conduct based on R^* . This procedure is repeated for 100 times for each tracker. Afterwards, the mean \bar{T} and standard deviation T_{std} of TSP values are calculated for each frame. Finally, we plot the *TSP band*, which is a band changing along with frame index and covers the range $[\bar{T} - T_{std}, \bar{T} + T_{std}]$, for every visual tracker.

The new experiment is carried out on video sequence *pets2000_c1* and the *TSP bands* are demonstrated in Figure 7. An ideal *TSP band* should be with small variance and centered around a relatively high mean. We can see that in Figure 7, RTCST and KMS tracker show similar variance but RTCST has a higher TSP mean. PF tracker illustrates smaller variance but suffers from very low accuracy. RTCST-B comes with the highest average TSP value while still achieves smallest standard deviation. The experiment result exhibits the unstable nature of KMS tracker with respect to original target position. Meanwhile, it also confirms our conjecture about the presence of high robustness when background information is taken into consideration.

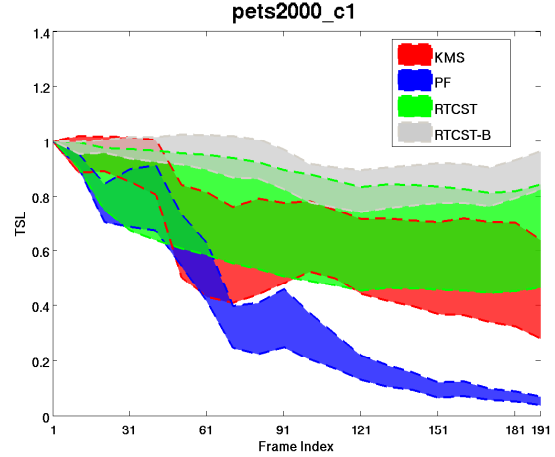


Fig. 7: Robustness Verification for visual trackers. The semi-transparent patches stand for the TSP bands of trackers. Note that here RTCST and RTCST-B are performed with D-100 features which is generated via random projection and 200 particles; PF tracker uses 500 particles.

VI. CONCLUSION AND FUTURE DIRECTIONS

In this paper, two enhanced CS-based visual tracking algorithms, namely, RTCST and RTCST-B are proposed. A customized OMP algorithm is designed to facilitate the proposed tracking algorithms. Hash kernel and random projection are employed to reduce the feature dimension of tracking application. In RTCST-B, a CS-based background model, which is termed CSBM, is utilized instead of noise templates. The new trackers achieves significantly higher efficiency compared with their prototype—the ℓ_1 tracker. The remarkable speed growth, which is up to 6271 times, makes CS-based visual trackers qualified for real-time applications. Meanwhile, our methods also obtain higher accuracy than off-the-shelf tracking algorithms, *i.e.*, PF tracker and KMS tracker. Particularly, RTCST-B achieves consistently highest accuracy and robustness thanks to the exploitation of background information. In short words, the proposed RTCST and RTCST-B are sufficiently fast for real-time visual tracking and more accurate and robust than conventional trackers.

For future topics, we believe that one low-hanging fruit is employing the trick mentioned in [22] by Tropp *et al.* to accelerate the OMP procedure furthermore. Another promising direction is to take color information into consideration because in many scenarios, color-based classification is more discriminant than the intensity-based one. The third direction of future research is treating different part of the target, *e.g.* left-top quarter and middle-bottom quarter, as different classes. As a result, a multiple classification is conduct within CS framework. The obtained likelihood for each particle then becomes a vector comprised of the confidences associated with various target parts. Because the time consumptions for binary and multiple classification are the same when using CS-based manner, we actually obtain more information at the same cost. If we can find a reasonable way to exploit the extra information for tracking, more accurate and robust result is likely to be obtained.

TABLE IV: Running time of visual trackers for one frame (ms). Note that every time consumption based on Matlab implementation is labeled by signal “*”. The notations of algorithm names are the same to those used in Table III.

		cubicle	dp	car4	car11	fish	pets2000_c1	pets2001_c1	pets2002_p1	pets2004_p1	pets2004-2_p1
KMS		22 ± 0	17 ± 0	60 ± 0	15 ± 0	36 ± 0	83 ± 0	40 ± 0	22 ± 0	21 ± 0	31 ± 0
PF	PN100	18 ± 0	17 ± 0	173 ± 0	22 ± 0	39 ± 0	199 ± 0	35 ± 0	28 ± 0	53 ± 0	381 ± 0
	PN200	27 ± 0	20 ± 0	321 ± 0	32 ± 0	56 ± 0	279 ± 0	37 ± 0	44 ± 0	82 ± 0	734 ± 0
	PN500	40 ± 0	37 ± 0	770 ± 0	65 ± 0	139 ± 0	631 ± 0	45 ± 0	83 ± 0	184 ± 0	1727 ± 0
R-D25-Rand	PN100	84 ± 3*	100 ± 4*	115 ± 2*	103 ± 4*	114 ± 4*	105 ± 3*	103 ± 4*	131 ± 3*	109 ± 2*	117 ± 3*
	PN200	148 ± 9*	152 ± 11*	193 ± 5*	186 ± 11*	198 ± 15*	186 ± 9*	177 ± 11*	223 ± 8*	168 ± 7*	198 ± 5*
R-D50-Rand	PN100	155 ± 4*	171 ± 4*	189 ± 4*	197 ± 5*	168 ± 10*	192 ± 7*	187 ± 5*	188 ± 5*	169 ± 5*	167 ± 4*
	PN200	276 ± 18*	301 ± 19*	337 ± 13*	358 ± 20*	333 ± 33*	334 ± 23*	334 ± 24*	347 ± 15*	286 ± 13*	336 ± 13*
R-D100-Rand	PN100	477 ± 21*	474 ± 17*	480 ± 23*	535 ± 10*	435 ± 32*	473 ± 11*	496 ± 26*	478 ± 18*	439 ± 17*	481 ± 13*
	PN200	825 ± 97*	742 ± 101*	939 ± 21*	968 ± 71*	870 ± 101*	798 ± 86*	863 ± 94*	863 ± 42*	681 ± 58*	872 ± 29*
R-D25-Hash	PN100	91 ± 3*	92 ± 4*	109 ± 3*	109 ± 3*	103 ± 5*	110 ± 4*	108 ± 4*	131 ± 3*	102 ± 3*	108 ± 2*
	PN200	166 ± 9*	161 ± 7*	172 ± 7*	191 ± 13*	193 ± 14*	204 ± 9*	195 ± 12*	217 ± 10*	161 ± 9*	194 ± 6*
R-D50-Hash	PN100	57 ± 1*	54 ± 1*	67 ± 1*	62 ± 2*	56 ± 1*	70 ± 1*	65 ± 1*	70 ± 1*	59 ± 2*	59 ± 1*
	PN200	96 ± 2*	96 ± 2*	118 ± 1*	114 ± 2*	101 ± 2*	118 ± 2*	114 ± 3*	116 ± 2*	109 ± 3*	108 ± 1*
R-D100-Hash	PN100	73 ± 1*	85 ± 2*	87 ± 2*	86 ± 1*	84 ± 3*	96 ± 3*	83 ± 2*	96 ± 5*	78 ± 2*	82 ± 1*
	PN200	138 ± 4*	154 ± 4*	148 ± 3*	156 ± 2*	157 ± 2*	162 ± 4*	146 ± 4*	169 ± 3*	134 ± 2*	159 ± 1*
RB-D25-Rand	PN100	—	—	—	—	—	167 ± 5*	175 ± 6*	204 ± 6*	142 ± 23*	184 ± 4*
	PN200	—	—	—	—	—	305 ± 11*	330 ± 19*	316 ± 26*	237 ± 32*	331 ± 18*
RB-D50-Rand	PN100	—	—	—	—	—	187 ± 7*	228 ± 4*	222 ± 7*	157 ± 33*	211 ± 4*
	PN200	—	—	—	—	—	389 ± 27*	500 ± 29*	397 ± 28*	295 ± 74*	427 ± 21*
RB-D100-Rand	PN100	—	—	—	—	—	215 ± 4*	246 ± 3*	248 ± 7*	148 ± 36*	253 ± 8*
	PN200	—	—	—	—	—	456 ± 17*	534 ± 23*	438 ± 38*	318 ± 75*	461 ± 45*
RB-D25-Hash	PN100	—	—	—	—	—	162 ± 7*	177 ± 8*	180 ± 11*	131 ± 27*	178 ± 8*
	PN200	—	—	—	—	—	274 ± 18*	377 ± 28*	306 ± 29*	227 ± 41*	351 ± 15*
RB-D50-Hash	PN100	—	—	—	—	—	95 ± 2*	88 ± 3*	106 ± 2*	85 ± 2*	94 ± 1*
	PN200	—	—	—	—	—	174 ± 5*	166 ± 2*	176 ± 3*	154 ± 8*	174 ± 3*
RB-D100-Hash	PN100	—	—	—	—	—	121 ± 2*	106 ± 2*	127 ± 3*	111 ± 3*	114 ± 2*
	PN200	—	—	—	—	—	220 ± 6*	211 ± 7*	229 ± 5*	207 ± 8*	217 ± 5*
LIT-Matlab		2.7e5 ± 1255*	8.7e4 ± 1660*	—	1.8e5 ± 2402*	—	—	1.6e5 ± 1944*	—	3.7e5 ± 1857*	—
LIT-C++		3.2e4 ± 506	1.4e4 ± 320	—	3.8e4 ± 1417	—	—	3.4e4 ± 484	—	1.02e5 ± 607	—

REFERENCES

- [1] X. Mei and H. Ling, “Robust visual tracking using ℓ_1 minimization,” in *Proc. IEEE Int. Conf. Comp. Vis.*, Kyoto, Japan, 2009, pp. 1436–1443.
- [2] T. F. Cootes, G. J. Edwards, and C. J. Taylor, “Active appearance models,” *IEEE Trans. Pattern Anal. Mach. Intell.*, pp. 484–498, 1998.
- [3] D. Comaniciu, V. Ramesh, and P. Meer, “Kernel-based object tracking,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 25, pp. 564–577, 2003.
- [4] A. Yilmaz and M. Shah, “Contour-based object tracking with occlusion handling in video acquired using mobile cameras,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 26, pp. 1531–1536, 2004.
- [5] S. Avidan, “Support vector tracking,” *IEEE Trans. Pattern Anal. Mach. Intell.*, pp. 184–191, 2001.
- [6] D. Serby and L. V. Gool, “Probabilistic object tracking using multiple features,” in *Proc. IEEE Int. Conf. Patt. Recogn.*, 2004, pp. 184–187.
- [7] C. Shen, J. Kim, and H. Wang, “Generalized kernel-based visual tracking,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 20, pp. 119–130, 2010.
- [8] A. Doucet, S. Godsill, and C. Andrieu, “On sequential monte carlo sampling methods for bayesian filtering,” *Statistics and Computing*, vol. 10, no. 3, pp. 197–208, 2000.
- [9] C. Shen, M. J. Brooks, and A. van den Hengel, “Fast global kernel density mode seeking: applications to localization and tracking,” *IEEE Trans. Image Process.*, vol. 16, no. 5, pp. 1457–1469, 2007.
- [10] M. L. Cascia, S. Sclaroff, and V. Athitsos, “Fast, reliable head tracking under varying illumination: An approach based on registration of texture-mapped 3d models,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 22, pp. 322–336, 2000.
- [11] K. Shafique and M. Shah, “A non-iterative greedy algorithm for multi-frame point correspondence,” in *IEEE Trans. Pattern Anal. Mach. Intell.*, 2003, pp. 51–65.
- [12] O. Williams, A. Blake, and R. Cipolla, “Sparse bayesian learning for efficient visual tracking,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 27, pp. 1292–1304, 2005.
- [13] M. S. Arulampalam, S. Maskell, and N. Gordon, “A tutorial on particle filters for online nonlinear/non-gaussian bayesian tracking,” *IEEE Trans. Signal Process.*, vol. 50, pp. 174–188, 2002.
- [14] D. Comaniciu, V. Ramesh, and P. Meer, “Kernel-based object tracking,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 25, pp. 564–575, 2003.
- [15] F. Porikli, O. Tuzel, and P. Meer, “Covariance tracking using model update based on lie algebra,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2006, vol. 1, pp. 728–735.
- [16] S. Zhou, R. Chellappa, and B. Moghaddam, “Visual tracking and recognition using appearance-adaptive models in particle filters,” *IEEE Trans. Image Process.*, vol. 13, pp. 1434–1456, 2004.
- [17] J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry, and Y. Ma, “Robust face recognition via sparse representation,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 31, pp. 210–227, 2009.
- [18] Y. Tsaig and D. L. Donoho, “Compressed sensing,” *IEEE Trans. Inf. Theory*, vol. 52, pp. 1289–1306, 2006.
- [19] E. Candès, J. Romberg, and T. Tao, “Stable signal recovery from incomplete and inaccurate measurements,” *Communications on Pure and Applied Mathematics*, vol. 59, pp. 1207–1223, 2006.
- [20] V. Cevher, A. Sankaranarayanan, M. F. Duarte, D. Reddy, and R. G. Baraniuk, “Compressive sensing for background subtraction,” in *Proc. Eur. Conf. Comp. Vis.*, 2008, pp. 155–168.
- [21] Ali Cafer G., J. H. McClellan, J. Romberg, and W. R. Scott, “Compressive sensing of parameterized shapes in images,” in *Proc. IEEE Int. Conf. Acoust., Speech., Signal Process.*, 2008, pp. 1949–1952.
- [22] J. A. Tropp and A. C. Gilbert, “Signal recovery from random measurements via orthogonal matching pursuit,” *IEEE Trans. Inf. Theory*, vol. 53, pp. 4655–4666, 2007.
- [23] Q. Shi, J. Petterson, G. Dror, J. Langford, A. Smola, A. Strehl, and S. V. N. Vishwanathan, “Hash kernels,” in *Proc. Int. Workshop Artificial Intell. & Statistics*, 2009.
- [24] S. Mallat and Z. Zhang, “Matching pursuit with time-frequency dictionaries,” *IEEE Trans. Signal Process.*, vol. 41, pp. 3397–3415, 1993.
- [25] Y. C. Pati, R. Rezaiifar, Y. C. Pati R. Rezaiifar, and P. S. Krishnaprasad, “Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition,” in *Proceedings of the 27th Annual Asilomar Conference on Signals, Systems, and Computers*, 1993, pp. 40–44.
- [26] G. Davis, S. Mallat, and Z. Zhang, “Adaptive time-frequency decompositions with matching pursuits,” *Optical Engineering*, vol. 33, 1994.
- [27] C. Stauffer and W. E. L. Grimson, “Adaptive background mixture models for real-time tracking,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 1999, vol. 2, pp. 246–252.
- [28] M. Isard and J. Maccormick, “Bramble: A bayesian multiple-blob tracker,” in *Proc. IEEE Int. Conf. Comp. Vis.*, 2001, vol. 2, pp. 34–41.
- [29] T. Zhao, R. Nevatia, and F. Lv, “Segmentation and tracking of multiple humans in complex situations,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 30, pp. 1198–1211, 2001.
- [30] A.S. MOSEK, “The MOSEK optimization software,” *Online at <http://www.mosek.com>*, 2010.
- [31] M. Everingham, L. V. Gool, C.K.I. Williams, J. Winn, and A. Zisserman, “The PASCAL visual object classes (VOC) challenge,” *Int. J. Comp. Vis.*, vol. 88, no. 2, pp. 303–338, 2010.

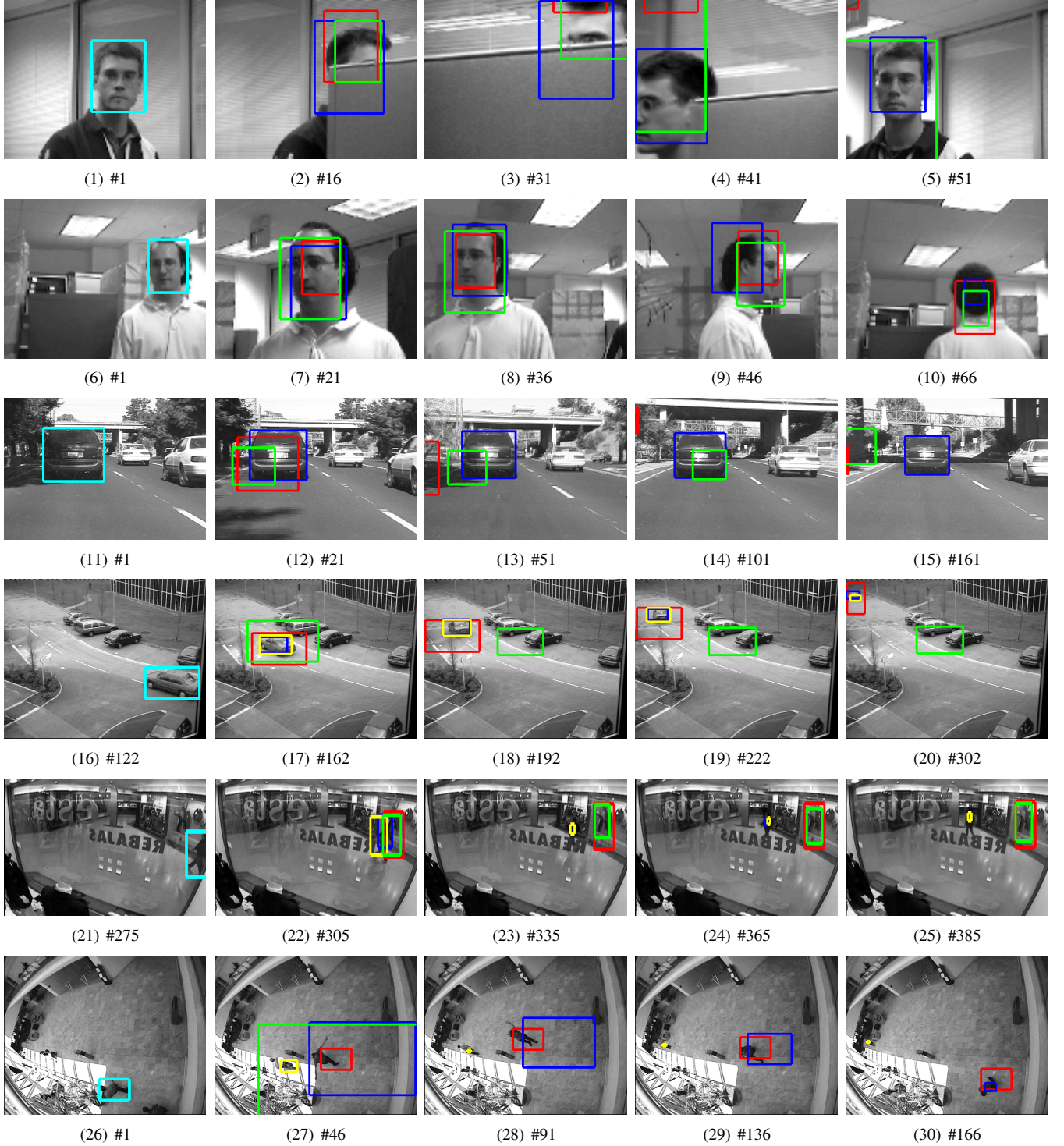


Fig. 8: Tracking results shown as rectangles for 6 video sequences, namely, *cubicle*, *dp*, *car4*, *pets2000_c1*, *pets2002_p1* and *pets2004_p1*. Symbol # x stands for the x th frame. The initial target position is shown in light blue while the red, green, dark blue and yellow rectangle denote the tracked area by KMS tracker, PF tracker (PN500), RTCST (D100-Rand-PN200) and RTCST-B (D100-Rand-PN200) respectively. For a certain tracker, the illustrated result is the one with the highest TSP value among all the associated results. PF tracker extends the tracking region to the whole scene in the latter frames on *pets2004_p1*, this is why we can not see the green rectangle in these frames. RTCST and RTCST-B tracking the similar regions for the last frame on *pets2002_p1* and the yellow rectangle covers the blue one.

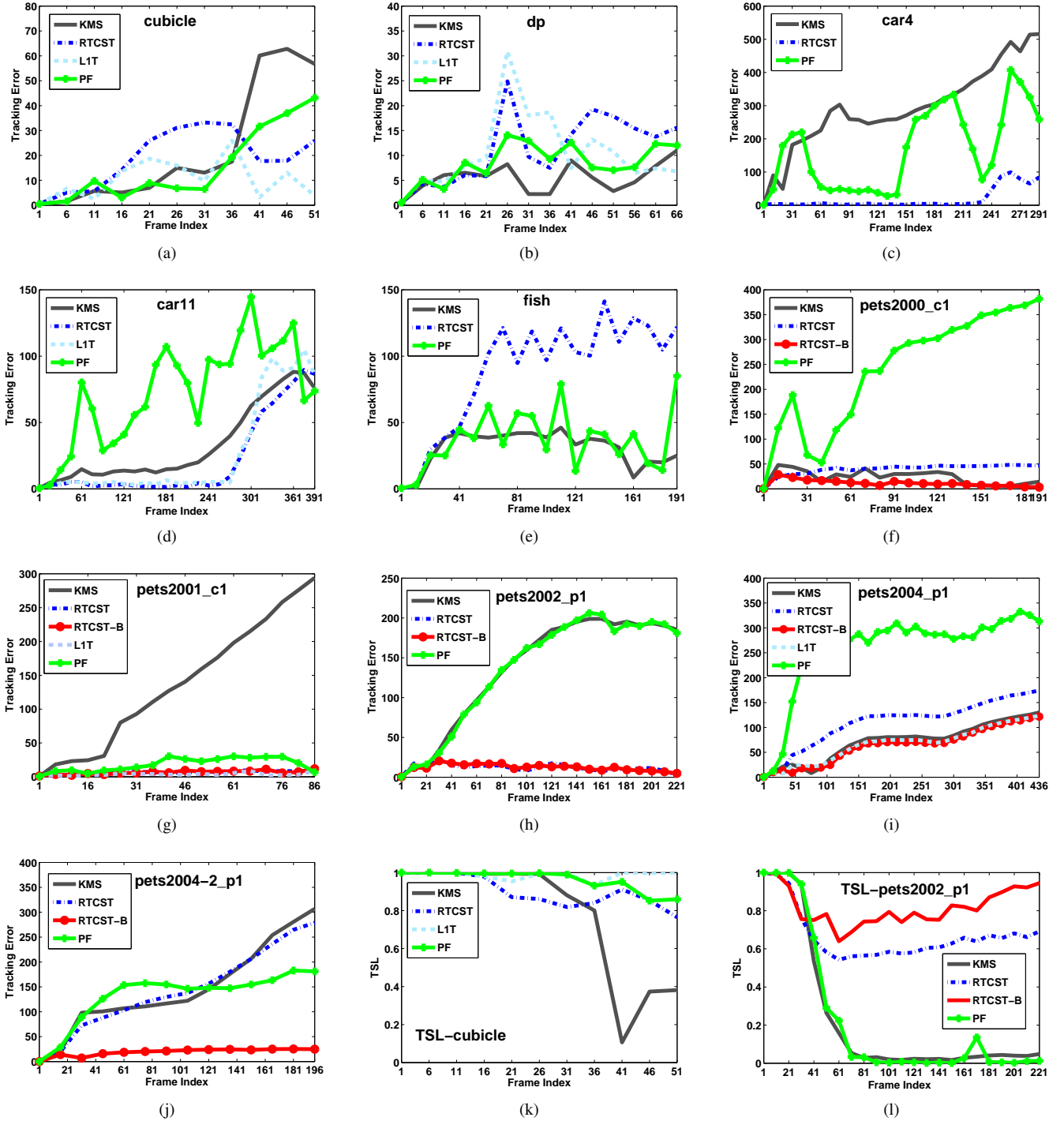


Fig. 9: The tracking errors and TSP values changing along with the frame index. All the visual trackers employ the optimal parameters, *i.e.*, 500 particles for PF tracker; 200 particles and Dimension-100 for both RTCST and RTCST-B.